

Ameriflux Data Curation Portal Status and Learnings to Date

Berkeley Water Center Microsoft TCI
4/26/2006

Preamble

- Thanks for taking the time!
- This deck is a variant of an April status deck used internally in MSFT and with other e-science groups
- We're made some progress
- We're still in tool learning mode – many of the tools are new to us and our knowledge is the biggest stumbling block
- We're looking for feedback to improve and suggestions for next steps

2

Outline

- Ameriflux Data Overview (reminder)
- What does the "original" data look like?
- What schema are we using? Why?
- What portal workflows are we going to implement first?
- Current status – a few fun pictures
- Proposed next steps

3

Ameriflux Overview

149 Sites across the Americas

Each site reports a minimum of 22 common measurements.

Communal science – each principle investigator acts independently to prepare and publish data.

Data published to and archived at Oak Ridge.

<http://public.ornl.gov/ameriflux/>



Ameriflux Data Access Today

There are multiple (different) paths to the published data
Overview summary of (site:year) indicating data that has been submitted:

http://public.ornl.gov/ameriflux/viewstatus_Ameriflux.cfm

Single site .csv data download and simple graphics:

http://cdiac.esd.ornl.gov/programs/ameriflux/data_system/aa_mer.html

Multiple site .csv data download:

http://cdiac.esd.ornl.gov/programs/ameriflux/data_system/aa_lamer_pf.html

5

Ameriflux Data Access Today (fine print)

The single site page access includes only a subset of the sites indicated as having submitted data.

The multiple site download is a subset of the single site list.

A (hidden) ftp site holds older data (undocumented status) and on request data for specific investigators.

The available data is a mix of direct measurements and computation.

Data is subject to change at any time.

Data may differ from site to site. Some even include code!

Exact metrics differ from investigator to investigator. (Humidity may be .73 or 73.; Nulls and gaps may be null, zero, or -9999, or...)

Metadata is also dirty. (Eg reported latitude/longitude may vary across pages referencing site and within produced downloaded .csv file.)

The multi-site or multi-year downloads yield one big and likely unwieldy .csv file. (Can't import into excel, notepad, etc.) If the download actually completes.

Can view plots of one or two variables per site – little ability to compare sites or preview data available before download

6

What was wrong?

- All text fields should be unicode
 - » SQL Server Integration Services (SSIS) default is unicode; anything else requires more complex logic to convert
- No way of identifying a dataset across time
 - » Our "original" data is actually processed by investigators who can republish the "same" dataset.
- No built-in ability to derive a working dataset for scientific use
 - » Original data is gap filled, recalibrated, as well as used for graphs and statistics
 - » Ability to version is core to the analysis – is this calibration algorithm better/worse than another?
- Data vector structure leads to accumulating code that knows about that structure
 - » Not as reusable for other similar data or, more importantly, when federating two very different types of data

13

Lesson #2: Original data is in the eye of the beholder

One investigator's "raw" or "original" data is another's processed data.
So, after initial load of "original" data, derive a working dataset for scientific use.

Lesson #3: Abstracted datasets are your friends

Abstracted (normalized, pivoted) datasets should allow subsequent analysis tools to be common. In other words, interesting time series and statistical analyses can be written once and not per datum type.

Lesson #4: Keep the data load as simple as possible.

While the documentation explains how to extend SSIS with amazing scripting and complex flows, simple flows are easier to debug when the data is dirty. Even if you load into a dedicated database. Oh, and don't scrimp on space – keep everything. Also, add "superfluous" columns to track anything you might wonder about later.

Lesson #4 We need a clear data usage model

There are lots of twisty passages that computer scientist can go down. Which helps you more?

Datasets

- Datasets are collections of data used for a specific analysis or the result of measurement or simulation or library research or other information production activity
- All datasets have a unique identifier.
- Datasets may be *published* or *private*.
- Datasets are stored in an *original* or vectorized format or a *working* or atomized format.

18

Published Datasets Usage

■Published datasets:

- Are discoverable by and usable by other investigators, although some access policy may apply.
- Are curated in an archive
- Have contents that are fixed (unchangeable)

■All datasets used to derive a dataset must be published when that dataset is published.

- The unique identifiers are used to derive a published dataset are published with the published dataset.
- The unique identifiers can be used to establish provenance over time and subsequent use by other scientists.

■Published datasets may be:

- Published* by uploading to an archiving portal
- Downloaded* from an archiving or other caching portal
- <Open question> *Deleted* when there is no other published dataset that has been derived from the dataset.

19

Private Dataset Usage

■Private datasets:

- Are private to an individual investigator or small group of investigators.
- May change at any time.
- Before publishing, the dataset unique identifier can be used only by an individual scientist or small collaborating group of scientists.

■Private dataset operations:

- Import* is the ingest of an ad hoc shared dataset.
 - Analogous to download, but no discovery, access control or guarantee of stability.
- Export* is the emission of an ad hoc shared dataset.
 - Analogous to publishing, but no discovery, access control or guarantee of stability.
- Creation* can result from actual measurements, simulation, literature research or any other means.
- Derivation* is the result of some sort of data transformation (eg gap fill or calibration) on one or more existing published or private datasets.
- Appending* adds data to an existing unpublished dataset.
- Pruning* removes data from an existing unpublished dataset.
- Deletion* eliminates a non-published dataset.

20

Working or Atomized Datasets

■Normalized schema:

- Allows a wide variety of datum types to be mixed in a single dataset
- Enables analysis code and UI driving that code to be independent of datum type and private data formats.
- Is what an experienced database programmer would use!
- Not intuitive to most scientific programmers.

■Atomization requires:

- All values must be interpreted indirectly through datum identifiers.
- A specific datum identifier includes the name and units of the datum.
- The data format is rigidly determined by the datum identifier.
 - For example, the datum identifier used to report humidity as ".73" differs from that used to report humidity as ".73".
- A clearinghouse for datum identifiers is necessary. Several are evolving. We simply assume there is one.
- <Open question> how is a scientist going to be able to do that mapping simply and robustly given the richness?

21

Original or Vectorized Datasets

■Original datasets have private vectorized table format that can be mapped to a CSV file.

- CSV files are the lingua franca today for data interchange and simple analysis via Excel or other tools
- The column layout and data representation (eg units and scaling) are often private to investigators
- Collaborations often develop a common format for dataset publishing.
- Accommodating these vector formats are important to the scientist's workflow – especially since the data are often dirty and must be interpreted by the scientist.

■For each vector format, there is:

- A private data table format. That format must include a dataset unique identifier.
- Code that can import/export the vectorized data from/to a csv file. This includes a subset of the columns.
- Code that can convert the private table rows to/from atomized data.

22

New Ameriflux Database Schema



23

Working Dataset

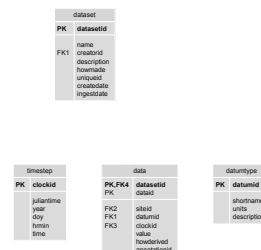
■Working datasets characterized by a name, creating investigator, description, and a GUID

■Actual data characterized by:

- Dataset**
- Site** (where collected)
- Datum** (for interpretation)
- Clock** (time)
- Value**
- How derived still TBD, but definitely includes "gap"

■Datum characterization eventually should correspond to GIS/HIS or other public schema

■(Recall that the Ameriflux data are published on regular time intervals)



24

Current Status

■ Visual Studio project leveraging SSIS for metadata and data loading of site, investigator, timestep and datum tables working from spreadsheets.

- Dynamic update of spreadsheet from web site not currently considered.

■ Raw data load project working from “Falk” spreadsheets.

- All "original data" loaded
- Working data set generated with air temperature (TA), CO2 flux (FC), CO2 in canopy (SC), and precipitation (PRECIP) for development
- Working data set is hourly only (temporary hack)

■ Met with potential contractor for portal download and data selection web service. Early sketches of UI developed jointly. Lawyers, lawyers, lawyers.

■ Talking with CUASHI HIS HODM team (David Tarboton, Reza Wahaji, Blair Jennings, Ilya Zaslavsky, et al) on common schema longer term.

31

Lesson #5: IT groups are like scientists.

The contractor's "customer model" for the code to be written should follow their practices when working with a medium-size business. The IT groups at such companies want to be able to debug and/or make minor changes in deployed applications without having to (re)call a contractor. That translates to being able to drag/drop or simple script or maybe even simple SQL, but no heavy programming. Note the similarities with how much a "scientist" wants to know about .NET, C#, etc.

Lesson #6: Our documentation sucks.

The tutorials (on line and in the documentation) are great for getting started, but don't give that many hints and/or assume a developer mind set or are evangelizing a technology. A non-expert (me) can browse that and learn all sorts of things but it's hit or miss.

For a good time, compare the online SSIS tutorial method for iterating over files to the SSIS book method. Is one better than the other? One is certainly *easier*.

Lesson #7:

(In progress) If you drag and drop to program, how do you debug?

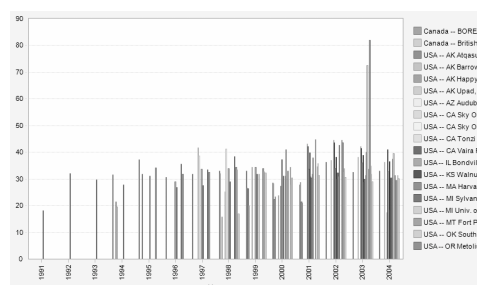
I've now managed to get breakpoints established and even a grid view going. But I've had best luck with a variant of the printf – adding a nvarchar(MAX) column to a table and stuffing suspicious text in along the way.

**Now for some fun:
What's wrong with the following
pictures?**

All of following slides were generated off the portal on my laptop.

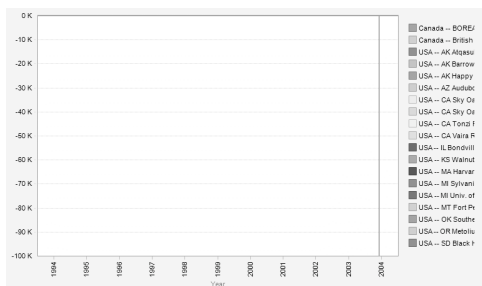
DISCLAIMER:
We're still very new to the tools and learning how to do things.

Maximum Air Temperature



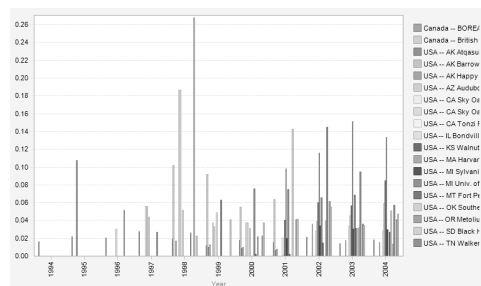
36

Minimum Precipitation



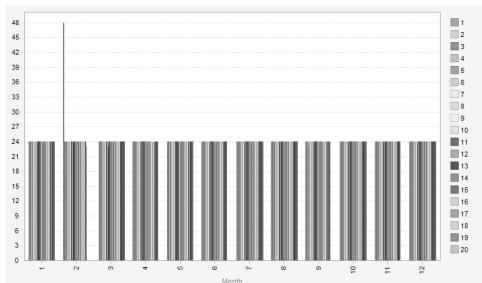
37

Mean Precipitation (w/o Vaira Ranch)



38

Count of Air Temperature measurements per month at Tonzi Ranch



39

Next Steps

■ Learn just enough MDX/SQL voodoo to enable your next questions (samples from Gretchen):

- Scatter plot temperature from Vaira and Tonzi Ranch
- How frequently did it rain in 2003 at Sky Oaks?
- How many days a year was the CO2 flux greater than average?
- How frequently did it rain in a given year?
- How many days in a given year had an above average carbon flux (fc) but a below average soil water content (swc)?
- When the temperature dropped more than 2 degrees in a one hour period (two consecutive measurements), what was the probability that rain would fall in the next two hours?
- ?????

■ Make the LBL portal available for your own data browsing.

■ Get the selection download UI working.

■ Continue discussions with the CUASHI HIS HODM folks.

40

Microsoft
Your potential. Our passion.™

© 2005 Microsoft Corporation. All rights reserved.
This presentation is for informational purposes only. Microsoft makes no warranties, express or implied, in this summary.